

The transposition distance for phylogenetic trees

Francesc Rosselló¹ and Gabriel Valiente²

¹ Department of Mathematics and Computer Science, Research Institute of Health Science (IUNICS), University of the Balearic Islands, E-07122 Palma de Mallorca, cesc.rossello@uib.es

² Algorithms, Bioinformatics, Complexity and Formal Methods Research Group, Department of Software, Technical University of Catalonia, E-08034 Barcelona, valiente@lsi.upc.edu

Abstract. The search for similarity and dissimilarity measures on phylogenetic trees has been motivated by the computation of consensus trees, the search by similarity in phylogenetic databases, and the assessment of clustering results in bioinformatics. The transposition distance for fully resolved phylogenetic trees is a recent addition to the extensive collection of available metrics for comparing phylogenetic trees. In this paper, we generalize the transposition distance from fully resolved to arbitrary phylogenetic trees, through a construction that involves an embedding of the set of phylogenetic trees with a fixed number of labeled leaves into a symmetric group and a generalization of Reidys-Stadler’s involution metric for RNA contact structures. We also present simple linear-time algorithms for computing it.

1 Introduction

The need for comparing phylogenetic trees arises when alternative phylogenies are obtained using different phylogenetic methods or different gene sequences for a given set of species. The comparison of phylogenetic trees is also essential to performing phylogenetic queries on databases of phylogenetic trees [8]. Further, the need for comparing phylogenetic trees also arises in the comparative analysis of clustering results obtained using different clustering methods or even different distance matrices, and there is a growing interest in the assessment of clustering results in bioinformatics [6].

A number of metrics for phylogenetic tree comparison are known, including the partition (or symmetric difference) metric [9,12], the nearest-neighbor interchange metric [19], the subtree transfer distance [1], the metric from the crossover method [13], the quartet metric [5], the metric from the nodal distance algorithm [3]. One of the simplest and easiest to compute metrics proposed so far, the transposition distance [17], is only defined for fully resolved trees. But phylogenetic analyses often produce phylogenies with polytomies, that is, phylogenetic trees that are not fully resolved. As a matter of fact, at the time of this writing, more than a 66.5% of the phylogenies contained in TreeBASE have polytomies.

In this paper, we generalize to arbitrary phylogenetic trees this transposition distance, through a new definition of it. This new distance is directly inspired on the one hand by the matching representation of phylogenetic trees [4,16] and on the other hand by the *involution metric* for RNA contact structures [11,14].

The matching representation $M(T)$ of a phylogenetic tree $T = (V, E)$ with n leaves labeled $1, \dots, n$ describes T injectively as a partition of $\{1, \dots, |V| - 1\}$. If T is fully

resolved, which is the particular case considered in [4], then all members of this partition are 2-elements sets, and then, since $|V| = 2n - 1$, it defines an undirected 1-regular graph $(\{1, \dots, 2n - 2\}, M(T))$. Reidys and Stadler defined the *involution metric* on 1-regular graphs, by associating to each such a graph the permutation given by the product of the transpositions corresponding to its edges, and then using the *canonical metric* in the symmetric group SS_{2n-2} (the least number of transpositions necessary to transform one permutation into another) to compare these permutations. The translation of this metric to matching representations yields twice the matching distance defined in [17]. Unfortunately, no meaningful generalization to arbitrary graphs of Reidys and Stadler’s metric is known, the main drawback being the difficulty of associating injectively a well-defined permutation to an arbitrary graph.

Now, if T is not fully resolved, the members of $M(T)$ are no longer pairs of numbers, and therefore they do not define a graph, at least not directly. Actually, the approach that we take in this paper can be understood as if we represented each member $\{i_1, \dots, i_k\}$ of $M(T)$, with $i_1 < \dots < i_k$, as a cyclic directed graph with arcs $(i_1, i_2), \dots, (i_{k-1}, i_k), (i_k, i_1)$, and $M(T)$ as the sum of these cyclic graphs. Now, generalizing Reidys-Stadler’s approach, we associate to every such a cyclic directed graph the cyclic permutation (i_1, \dots, i_k) (if $k = 2$, it is a transposition), and we describe $M(T)$ by means of the product of the cyclic permutations associated to its members: since these members are disjoint to each other, this product is well-defined. This defines an embedding of the set of phylogenetic trees with n leaves labeled $1, \dots, n$ into the symmetric group SS_{2n-2} . The transposition distance is obtained by translating the canonical metric on SS_{2n-2} into a distance for phylogenetic trees through this embedding. This transposition distance measures the least number of certain simple operations (splitting sets of children, joining sets of children, interchanging children) that are necessary to transform one tree into another, and it can be easily computed in linear time. Therefore it satisfies the requirements of “computational simplicity” and “good theoretical basis” that are required to any distance notion on phylogenetic trees [2].

2 Matching Representation of Phylogenetic Trees

Throughout this paper, by a *phylogenetic tree* we mean a *rooted tree with injectively labeled leaves and without outdegree 1 nodes*. Thus, a phylogenetic tree is a directed finite graph $T = (V, E)$ containing a distinguished node $r \in V$, called the *root*, such that for every other node $v \in V$ there exists one, and only one, path from the root r to v . The *children* of a node v in a tree $T = (V, E)$ are those nodes $w \in V$ such that $(v, w) \in E$. The *outdegree* of a node is the number of its children. The nodes without children are the *leaves* of the tree, and the remaining nodes are called *internal*: since we assume that no node has outdegree 1, every internal node has at least 2 children. The set of leaves of T is denoted by $\mathcal{L}(T)$. The *height* of a node v in a tree T is the length of a longest directed path from v to a leaf. Thus, the nodes with height 0 are the leaves, the nodes with height 1 are the nodes all whose children are leaves, and so on.

The leaves of a phylogenetic tree are injectively labeled in a fixed, but arbitrary, ordered set: these labels are called *taxa*. In practice, if the tree has n leaves, we shall identify their labels with $1, \dots, n$, ordered in the usual increasing way. The label associated to a leaf $v \in V$ will be denoted by $\ell(v)$.

We shall denote by \mathcal{T}_n the set of all phylogenetic trees with n leaves labeled $1, \dots, n$ (up to label-preserving isomorphisms of rooted trees).

Definition 1. *The bottom-up ordering (cf. [4, 18]) of a phylogenetic tree $T = (V, E) \in \mathcal{T}_n$ is the injective mapping*

$$\ell : V \rightarrow \{1, \dots, |V|\}$$

defined by the following properties:

- (a) *If $v \in \mathcal{L}(T)$, then $\ell(v)$ is its label.*
- (b) *If $\text{height}(u) < \text{height}(v)$, then $\ell(u) < \ell(v)$.*
- (c) *If $0 < \text{height}(u) = \text{height}(v)$ and*

$$\min\{\ell(x) \mid x \in \text{children}(u)\} < \min\{\ell(x) \mid x \in \text{children}(v)\},$$

then $\ell(u) < \ell(v)$.

It is straightforward to notice that this bottom-up ordering is unique, and it can be computed in time linear in the size of the tree by bottom-up tree traversal techniques [18]. First, the leaves of T are labeled by their label in $\{1, \dots, n\}$. Then, the height 1 nodes are labeled from $n + 1$ on in the order given by the smallest label of their children: i.e., the height 1 node with the smallest child label is assigned label $n + 1$, the height 1 node with the next-smallest child label is assigned label $n + 2$, etc. And this procedure is continued for consecutively increasing heights. The detailed pseudocode is given in Algorithm 1.

Example 1. Fig. 1 shows the Tree T166c11x6x95c08c56c38 in TreeBASE and its bottom-up ordering after sorting its taxa alphabetically.

The next definition generalizes the perfect matching representation of binary, or fully resolved, trees [4, 16].

Definition 2. *Let $T = (V, E)$ be a phylogenetic tree with n leaves labeled $1, \dots, n$, and let $\ell : V \rightarrow \{1, \dots, |V|\}$ be its bottom-up ordering. The matching representation $M(T)$ of T is the partition of $\{1, \dots, |V| - 1\}$ defined as follows:*

$$M(T) = \{\ell(\text{children}(u)) \mid u \in V - \mathcal{L}(T)\}.$$

Example 2. The matching representation of the tree in Fig. 1 is the partition of $\{1, \dots, 14\}$ given by

$$\left\{ \{1, 5, 7, 9\}, \{4, 6, 10\}, \{2, 11\}, \{8, 13\}, \{3, 12, 14\} \right\}.$$

```

begin
  foreach node  $v$  of  $T$  do
    if  $v$  is a leaf node of  $T$  then
      set  $\ell(v)$  to the index of  $\ell(v)$  in  $L$ 
    else
       $\ell(v) := 0$ 
  end
   $i := |L|$ 
  foreach level  $h$  of  $T$  from the leaves up to the root do
    let  $S$  be the set of nodes of  $T$  at level  $h$ , ordered by label
    foreach  $v \in S$  do
      let  $w$  be the parent of  $v$  in  $T$ 
      if  $\ell(w) = 0$  and  $\text{height}(w) = h + 1$  then
         $i := i + 1$ 
         $\ell(w) := i$ 
      end
    end
  end
  return  $M$ 
end

```

Algorithm 1: Bottom-up ordering. Given an ordered set L and a phylogenetic tree T with leaves bijectively labeled in L , the algorithm computes the bottom-up ordering of T .

It is clear that, once the bottom-up ordering of T has been obtained, the set $M(T)$ can be produced in linear time in the size of the tree. Furthermore, the following two results are straightforward.

Corollary 1. *For every $T = (V, E) \in \mathcal{T}_n$, $|M(T)| = |V| - n$.*

Corollary 2. *For every $T_1, T_2 \in \mathcal{T}_n$, if $M(T_1) = M(T_2)$, then $T_1 = T_2$.*

3 The transposition distance

For every $m \geq 1$, let SS_m denote the symmetric group of permutations of $\{1, \dots, m\}$. By a *cycle* in SS_m we understand a cyclic permutation $(i_1, i_2, \dots, i_k) \in SS_m$, with $k \geq 2$, that sends i_1 to i_2 , i_2 to i_3 , ..., i_{k-1} to i_k , and i_k to i_1 , leaving fixed the remaining elements of $\{1, \dots, m\}$. Recall that the inverse of a cycle (i_1, i_2, \dots, i_k) is $(i_1, i_2, \dots, i_k)^{-1} = (i_k, i_{k-1}, \dots, i_1)$: the permutation that sends i_k to i_{k-1} , i_{k-1} to i_{k-2} , ..., i_2 to i_1 , and i_1 to i_k . The *length* of a cycle (i_1, i_2, \dots, i_k) is the number k of elements it moves.

The *cycle associated to a subset* $S = \{i_1, \dots, i_k\}$, with $i_1 < \dots < i_k$ and $k \geq 2$, of $\{1, \dots, m\}$, is $\kappa(S) := (i_1, i_2, \dots, i_k) \in SS_m$. If $k = 1$, i.e., if S is a singleton, then $\kappa(S)$ is the identity in SS_m , which we do not consider a cycle.

Definition 3. *The matching permutation $\pi(T)$ associated to a phylogenetic tree $T = (V, E) \in \mathcal{T}_n$ is the permutation of $\{1, \dots, |V| - 1\}$ defined by the product of the sorted cycles associated to the members of its matching representation:*

$$\pi(T) = \prod_{u \in V - \mathcal{L}(T)} \kappa(\ell(\text{children}(u))).$$

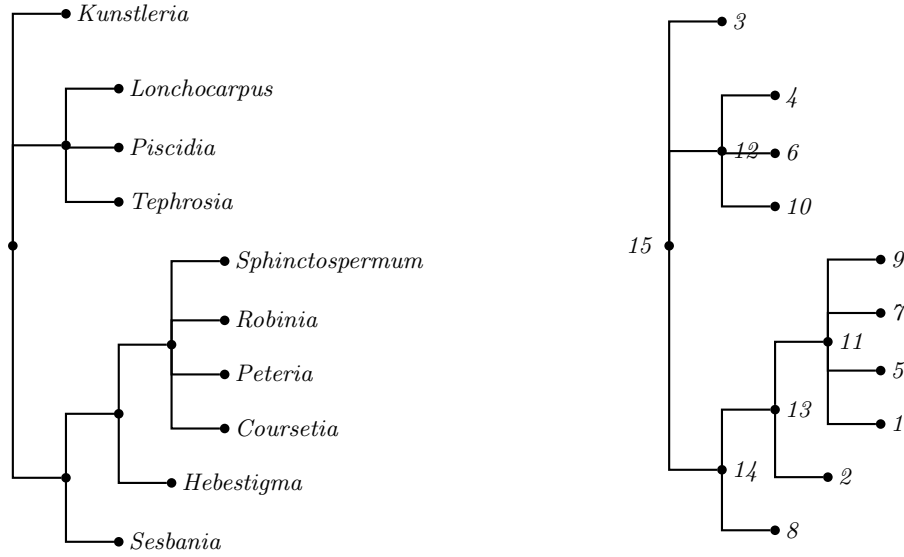


Fig. 1. A phylogenetic tree (left) and its bottom-up ordering (right).

Example 3. The matching permutation associated to the tree in Fig. 1 is the product of cycles

$$(1, 5, 7, 9)(4, 6, 10)(2, 11)(8, 13)(3, 12, 14) \in SS_{14},$$

i.e., the permutation

$$\begin{pmatrix} 1 & 2 & 3 & 4 & 5 & 6 & 7 & 8 & 9 & 10 & 11 & 12 & 13 & 14 \\ 5 & 11 & 12 & 6 & 7 & 10 & 9 & 13 & 1 & 4 & 2 & 14 & 8 & 3 \end{pmatrix}$$

If $u, v \in V - \mathcal{L}(T)$ are two different internal nodes of T , then $\ell(\text{children}(u)) \cap \ell(\text{children}(v)) = \emptyset$. Therefore, all cycles $\kappa(\ell(\text{children}(u)))$ appearing in the product defining $\pi(T)$ are disjoint to each other, and hence they commute with each other, which implies that this product is well defined.

Notice that no element in $\{1, \dots, |V| - 1\}$ remains fixed by $\pi(T)$, because every $\ell(\text{children}(u))$, with u internal, has at least two elements and every element in $\{1, \dots, |V| - 1\}$ is the bottom-up ordering label of a child of some internal node. Now, if $T = (V, E)$ is a phylogenetic tree with n leaves, then $|V| \leq 2n - 1$, the equality holding if and only if T is binary. To be able to compare matching permutations of phylogenetic trees with the same number of leaves n but different numbers of internal nodes, we shall understand henceforth that the matching permutation $\pi(T)$ belongs to SS_{2n-2} , leaving fixed the elements $|V|, \dots, 2n - 2$.

The following result is a direct consequence of the facts that the matching representation of a phylogenetic tree uniquely determines it and every permutation has a unique decomposition as a product of disjoint cycles of length ≥ 2 .

Proposition 1. *For every $T_1, T_2 \in \mathcal{T}_n$, if $\pi(T_1) = \pi(T_2)$, then $T_1 = T_2$.*

Remark 1. If we allow the existence of outdegree 1 nodes in our phylogenetic trees, then the last proposition is no longer true. Indeed, consider the trees in Fig. 2. The left-hand side one has matching representation $\{\{1, 2, 3\}, \{4\}\}$, while the right-hand side one has matching representation $\{\{1, 2, 3\}\}$. Therefore the matching permutation associated to both trees is $(1, 2, 3)$ (considered as an element of SS_4).



Fig. 2. Two trees with the same matching permutation.

Arguing as in [11, Cor. 1], we have the following result.

Theorem 1. *The mapping that associates to every pair (T_1, T_2) of phylogenetic trees with n leaves labeled in $\{1, \dots, n\}$, the least number $TD'(T_1, T_2)$ of transpositions necessary to represent the permutation $\pi(T_2)^{-1}\pi(T_1) \in SS_{2n-2}$, is a metric on \mathcal{T}_n .*

Proof. By Proposition 1, the mapping $\pi : \mathcal{T}_n \rightarrow SS_{2n-2}$ that sends every $T \in \mathcal{T}_n$ to its matching permutation $\pi(T)$ is an embedding. Then, since the mapping

$$d_{trans} : SS_{2n-2} \times SS_{2n-2} \rightarrow \mathbb{N}$$

defined by

$$d_{trans}(\pi_1, \pi_2) = \text{the least number of transpositions necessary to represent } \pi_2^{-1} \cdot \pi_1$$

is a metric on SS_{2n-2} (see, for instance, [11, Thm. 2]), the mapping

$$TD' : \mathcal{T}_n \times \mathcal{T}_n \rightarrow \mathbb{N} \\ (T_1, T_2) \mapsto d_{trans}(\pi(T_1), \pi(T_2))$$

is a metric on \mathcal{T}_n . □

Remark 2. Recall that the least number of transpositions required to represent a cycle of length k is $k - 1$, for instance through

$$(i_1, \dots, i_k) = (i_1, i_2)(i_2, i_3) \cdots (i_{k-1}, i_k),$$

and that the least number of transpositions required to represent a product of disjoint cycles is the sum of the least numbers of transpositions each cycle decomposes into, and hence the sum of the cycles' lengths minus the number of cycles.

The metric TD' satisfies the following property.

Proposition 2. *For every $T_1, T_2 \in \mathcal{T}_n$, $TD'(T_1, T_2)$ is an even integer smaller than $2n-2$.*

Proof. If $T_1, T_2 \in \mathcal{T}_n$ have m_1 and m_2 internal nodes, respectively, then each $\pi(T_i)$ ($i = 1, 2$) decomposes into m_i disjoint cycles: say $\pi(T_i) = C_{i,1} \cdots C_{i,m_i}$, with $C_{i,j}$ of length $k_{i,j}$. Then, by Remark 2, $\pi(T_i)$ has a decomposition into

$$\sum_{j=1}^{m_i} (k_{i,j} - 1) = \sum_{j=1}^{m_i} k_{i,j} - m_i = n + m_i - 1 - m_i = n - 1$$

transpositions. But then $\pi(T_2)^{-1}\pi(T_1)$ admits a decomposition into $2(n-1)$ transpositions. This entails that *every* decomposition of this permutation into a product of transpositions must involve an even number of them, and therefore that $TD'(T_1, T_2)$ is an even integer.

As far as the stated upper bound for $TD'(T_1, T_2)$ goes, notice that $\pi(T_2)^{-1}\pi(T_1)$ moves at most $2n-2$ elements and that if it is not the identity, then its decomposition into disjoint cycles has at least 1 cycle. Therefore, again by Remark 2, a minimal decomposition of this permutation into transpositions will involve at most $(2n-2) - 1$ transpositions, and since this number is even, this implies that $TD'(T_1, T_2) \leq 2n - 4$.

In other words, TD' is “artificially” multiplied by 2. Thus, we define a new metric on \mathcal{T}_n by dividing TD' by 2.

Definition 4. *The transposition distance on \mathcal{T}_n is*

$$\begin{aligned} TD : \mathcal{T}_n \times \mathcal{T}_n &\rightarrow \mathbb{N} \\ (T_1, T_2) &\mapsto \frac{1}{2} TD'(T_1, T_2) \end{aligned}$$

In this way, TD takes values in $\{0, 1, 2, \dots, n-2\}$.

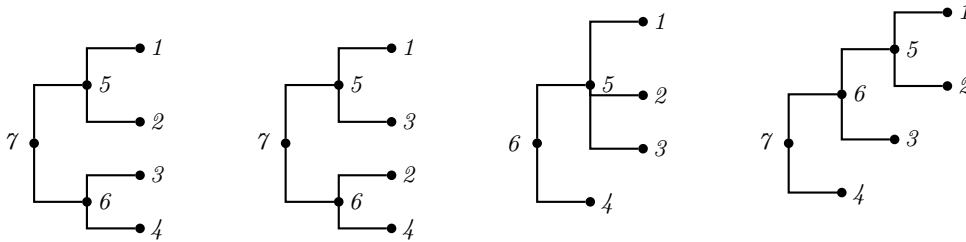


Fig. 3. From left to right, the phylogenetic trees T_1 , T_2 , T_3 , and T_4 in Example 4.

Table 1. Transposition distances between pairs of trees T_1, \dots, T_4 .

TD	T_1	T_2	T_3	T_4
T_1	0	1	2	1
T_2	1	0	2	2
T_3	2	2	0	2
T_4	1	2	2	0

Example 4. Let T_1, T_2, T_3, T_4 be the phylogenetic trees displayed in Fig. 3 (which we already give bottom-up ordered). Their matching permutations are

$$\begin{aligned}\pi(T_1) &= (1, 2)(3, 4)(5, 6), & \pi(T_2) &= (1, 3)(2, 4)(5, 6), \\ \pi(T_3) &= (1, 2, 3)(4, 5), & \pi(T_4) &= (1, 2)(3, 5)(4, 6)\end{aligned}$$

(understood as permutations in SS_6), and then

$$\begin{aligned}\pi(T_2)^{-1}\pi(T_1) &= (3, 1)(4, 2)(6, 5)(1, 2)(3, 4)(5, 6) = (1, 4)(2, 3) \\ \pi(T_3)^{-1}\pi(T_1) &= (3, 2, 1)(5, 4)(1, 2)(3, 4)(5, 6) = (2, 3, 5, 6, 4) \\ \pi(T_4)^{-1}\pi(T_1) &= (2, 1)(5, 3)(6, 4)(1, 2)(3, 4)(5, 6) = (3, 6)(4, 5) \\ \pi(T_3)^{-1}\pi(T_2) &= (3, 2, 1)(5, 4)(1, 3)(2, 4)(5, 6) = (1, 2, 5, 6, 4) \\ \pi(T_4)^{-1}\pi(T_2) &= (2, 1)(5, 3)(6, 4)(1, 3)(2, 4)(5, 6) = (1, 5, 4)(3, 2, 6) \\ \pi(T_4)^{-1}\pi(T_3) &= (2, 1)(5, 3)(6, 4)(1, 2, 3)(4, 5) = (2, 5, 6, 4, 3)\end{aligned}$$

which yields the distances between these trees given in Table 1.

The transposition distance between two phylogenetic trees can be easily computed in linear time. To prove it, we move to the more general setting of permutations and the graphs associated to them.

For every permutation $\pi \in SS_m$, the *directed graph* associated to π is the graph $G_\pi = (\{1, \dots, m\}, Q_\pi)$ with

$$Q_\pi = \{(i, j) \mid i \neq j \text{ and } \pi(i) = j\}.$$

The directed graph $G_{\pi^{-1}}$ associated to the inverse π^{-1} of a permutation π is obtained by reversing all arrows in G_π : thus, $Q_{\pi^{-1}} = Q_\pi^{-1}$ and $G_{\pi^{-1}} = G_\pi^{-1}$.

Given two permutations $\pi_1, \pi_2 \in SS_m$, by $G_{\pi_1} + G_{\pi_2}^{-1}$ we understand the 2-colored-arcs multigraph with set of nodes $\{1, \dots, m\}$, set of red arcs Q_{π_1} and set of blue arcs $Q_{\pi_2}^{-1}$. We shall say that a node of $G_{\pi_1} + G_{\pi_2}^{-1}$ is *unbalanced* when it is isolated in one, and only one, of the graphs $G_{\pi_1}, G_{\pi_2}^{-1}$ (which means that it is fixed by one, and only one, of the permutations π_1, π_2).

Proposition 3. *For every unbalanced node i of $G_{\pi_1} + G_{\pi_2}^{-1}$:*

- (1) *If i is isolated in G_{π_2} and $(i_0, i), (i, i_1) \in Q_{\pi_1}$ with $i_0 \neq i_1$, then replacing the red arcs (i_0, i) and (i, i_1) by a single red arc (i_0, i_1) increases $d_{trans}(\pi_1, \pi_2)$ by 1.*

- (2) If i is isolated in G_{π_2} and $(i, i_1), (i_1, i) \in Q_{\pi_1}$, removing the red arcs (i, i_1) and (i_1, i) increases $d_{trans}(\pi_1, \pi_2)$ by 1.
- (3) Similar properties hold if i is isolated in G_{π_1} but not in G_{π_2} and we modify the set of blue arcs.

Proof. (1) If $(i_0, i), (i, i_1) \in Q_{\pi_1}$, with $i_0 \neq i_1$, then $i_0 = \pi_1^{-1}(i)$ and $i_1 = \pi_1(i)$ and hence $(i, i_1)\pi_1(i_0) = i_1$, $(i, i_1)\pi_1(i) = i$, and $(i, i_1)\pi_1(j) = \pi_1(j)$ for every $j \neq i_0, i$. Therefore, replacing the arcs $(i_0, i), (i, i_1)$ by an arc (i_0, i_1) is equivalent to replacing π_1 by $(i, i_1)\pi_1$. So, it is enough to prove that, with the notations and assumptions of point (1),

$$d_{trans}(\pi_1, \pi_2) = d_{trans}((i, i_1)\pi_1, \pi_2) + 1.$$

To prove this equality, notice that, since i is fixed by π_2 , $\pi_2^{-1}\pi_1$ sends i_0 to i and i to $\pi_2^{-1}(i_1)$: let us denote this last index by j_1 .

If $j_1 = i_0$, then (i_0, i) is a cycle of $\pi_2^{-1}\pi_1$ and it appears in any decomposition of this permutation as a product of transpositions. But then both i and i_0 are fixed by $\pi_2^{-1}((i, i_1)\pi_1)$, and since $\pi_2^{-1}\pi_1$ and $\pi_2^{-1}((i, i_1)\pi_1)$ act exactly in the same way on the other elements, we deduce that

$$\pi_2^{-1}\pi_1 = (i_0, i)(\pi_2^{-1}((i, i_1)\pi_1))$$

and then $d_{trans}(\pi_1, \pi_2) = d_{trans}((i, i_1)\pi_1, \pi_2) + 1$ in this case.

If $j_1 \neq i_0$, then the cycle of $\pi_2^{-1}\pi_1$ moving i_0 has at least three elements:

$$(i_0, i, j_1, j_2, \dots, j_s), \quad \text{with } s \geq 1,$$

and thus it contributes $s + 1$ transpositions to a minimal decomposition of $\pi_2^{-1}\pi_1$ as a product of transpositions. Now, the cycle of $\pi_2^{-1}((i, i_1)\pi_1)$ that moves i_0 is

$$(i_0, j_1, j_2, \dots, j_s),$$

and it only contributes s transpositions to any decomposition of $\pi_2^{-1}((i, i_1)\pi_1)$ as a product of transpositions. Therefore, $d_{trans}(\pi_1, \pi_2) = d_{trans}((i, i_1)\pi_1, \pi_2) + 1$ also in this case.

- (2) If $(i, i_1), (i_1, i) \in Q_{\pi_1}$, then $\pi_1^{-1}(i) = \pi_1(i) = i_1$, and hence

$$(i, i_1)\pi_1(i) = i, \quad (i, i_1)\pi_1(i_1) = i_1,$$

and $(i, i_1)\pi_1(j) = \pi_1(j)$ for every $j \neq i, i_1$. Therefore, to remove the arcs $(i_1, i), (i, i_1)$ in this case means again to replace π_1 by $(i, i_1)\pi_1$. So, again in this case, it is enough to prove that, with the notations and assumptions of point (2),

$$d_{trans}(\pi_1, \pi_2) = d_{trans}((i, i_1)\pi_1, \pi_2) + 1.$$

Since i is fixed by π_2 , we have that $\pi_2^{-1}\pi_1$ sends i_1 to i and i to $\pi_2^{-1}(i_1)$: let us denote this last index by j_1 .

If $j_1 = i_1$, i.e., if i_1 is also fixed by π_2 , then (i, i_1) is a cycle of $\pi_2^{-1}\pi_1$ and it appears in any decomposition of this permutation as a product of transpositions. But then both i and i_1 are fixed by $\pi_2^{-1}((i, i_1)\pi_1)$ and

$$\pi_2^{-1}\pi_1 = (i_1, i)(\pi_2^{-1}((i, i_1)\pi_1))$$

and then $d_{trans}(\pi_1, \pi_2) = d_{trans}((i, i_1)\pi_1, \pi_2) + 1$.

If $j_1 \neq i_1$, then the cycle of $\pi_2^{-1}\pi_1$ moving i_1 has at least three elements:

$$(i_1, i, j_1, \dots, j_s), \quad \text{with } s \geq 1,$$

and thus it contributes $s+1$ transpositions to any decomposition of $\pi_2^{-1}\pi_1$ as a product of transpositions. Now, i is fixed by $\pi_2^{-1}((i, i_1)\pi_1)$ and the cycle of this permutation moving i_1 is

$$(i_1, j_1, j_2, \dots, j_s),$$

and it only contributes s transpositions to any decomposition of $\pi_2^{-1}((i, i_1)\pi_1)$ as a product of transpositions. Thus, again in this case, $d_{trans}(\pi_1, \pi_2) = d_{trans}((i, i_1)\pi_1, \pi_2) + 1$. \square

Proposition 4. *If $G_{\pi_1} + G_{\pi_2}^{-1}$ has no unbalanced node, then*

$$d_{trans}(\pi_1, \pi_2) = N(G_{\pi_1}, G_{\pi_2}^{-1}) - A(G_{\pi_1}, G_{\pi_2}^{-1}),$$

where $N(G_{\pi_1}, G_{\pi_2}^{-1})$ is the number of non-isolated nodes of $G_{\pi_1} + G_{\pi_2}^{-1}$ and $A(G_{\pi_1}, G_{\pi_2}^{-1})$ is the number of alternating cycles in $G_{\pi_1} + G_{\pi_2}^{-1}$, i.e., of cycles in this directed 2-colored-arcs multigraph such that two consecutive arcs have different colors.

Proof. If $G_{\pi_1} + G_{\pi_2}^{-1}$ has no unbalanced node, then every node either is isolated or has exactly one incoming and one outgoing arc of each color. This entails that $Q_{\pi_1} \sqcup Q_{\pi_2}$ decomposes into the union of arc-disjoint alternating cycles.

Now, every length $2k$ alternating cycle

$$(i_1, j_1), (j_1, i_2), (i_2, j_2), (j_2, i_3), \dots, (i_k, j_k), (j_k, i_1),$$

with $(i_\ell, j_\ell) \in Q_{\pi_1}$ for every $\ell = 1, \dots, k$ and $(j_\ell, i_{\ell+1}) \in Q_{\pi_2}$ for every $\ell = 1, \dots, k-1$ and $(j_k, i_1) \in Q_{\pi_2}$, corresponds to a length k cycle

$$(i_1, i_2, \dots, i_k)$$

of $\pi_2^{-1}\pi_1$ and hence it adds $k-1$ transpositions to any decomposition into transpositions of this permutation.

Therefore, if we denote by $\mathcal{A}(G_{\pi_1}, G_{\pi_2}^{-1})$ the set of alternating cycles in $G_{\pi_1} + G_{\pi_2}^{-1}$, we have that

$$\begin{aligned} d_{trans}(\pi_1, \pi_2) &= \sum_{C \in \mathcal{A}(G_{\pi_1}, G_{\pi_2}^{-1})} \left(\frac{\text{length}(C)}{2} - 1 \right) \\ &= \frac{1}{2} \sum_{C \in \mathcal{A}(G_{\pi_1}, G_{\pi_2}^{-1})} \text{length}(C) - |\mathcal{A}(G_{\pi_1}, G_{\pi_2}^{-1})| \\ &= \frac{1}{2} |Q_{\pi_1} \sqcup Q_{\pi_2}| - |\mathcal{A}(G_{\pi_1}, G_{\pi_2}^{-1})|. \end{aligned}$$

```

begin
  Compute the bottom-up orderings of  $T_1$  and  $T_2$ 
  Compute the matching representation  $M(T_1)$  and the directed graph
   $G_1 = (\{1, \dots, 2n-2\}, Q_1)$  associated to  $\pi(T_1)$ 
  Compute the matching representation  $M(T_2)$  and the directed graph
   $G_2 = (\{1, \dots, 2n-2\}, Q_2)$  associated to  $\pi(T_2)^{-1}$ 
   $d := 0$ 
   $N :=$  largest number appearing in  $M(T_1)$  or  $M(T_2)$ 
  while  $G_1 + G_2$  has unbalanced nodes do
    foreach angle  $\{(i_0, i), (i, i_1)\}$  in  $Q_1$  with  $i$  unbalanced and  $i_0 \neq i_1$  do
       $Q_1 := (Q_1 - \{(i_0, i), (i, i_1)\}) \cup \{(i_0, i_1)\}$ 
       $d := d + 1$ 
       $N := N - 1$ ;
    foreach angle  $\{(i_0, i), (i, i_1)\}$  in  $Q_2$  with  $i$  unbalanced and  $i_0 \neq i_1$  do
       $Q_2 := (Q_2 - \{(i_0, i), (i, i_1)\}) \cup \{(i_0, i_1)\}$ 
       $d := d + 1$ 
       $N := N - 1$ 
    foreach  $\{(i, i_1), (i_1, i)\}$  in  $Q_1$  with  $i$  unbalanced do
       $Q_1 := Q_1 - \{(i, i_1), (i_1, i)\}$ 
       $d := d + 1$ 
       $N := N - 1$  if  $i_1$  is not unbalanced,  $N := N - 2$  otherwise
    foreach  $\{(i, i_1), (i_1, i)\}$  in  $Q_2$  with  $i$  unbalanced do
       $Q_2 := Q_2 - \{(i, i_1), (i_1, i)\}$ 
       $d := d + 1$ 
       $N := N - 1$  if  $i_1$  is not unbalanced,  $N := N - 2$  otherwise
  Compute the number  $A$  of alternating cycles in the resulting directed multigraph  $G_1 + G_2$ , by
  traversing them
   $TD(T_1, T_2) := (d + N - A)/2$ 
end

```

Algorithm 2: Transposition distance. Given phylogenetic trees $T_1, T_2 \in \mathcal{T}_n$, the algorithm computes the transposition distance $TD(T_1, T_2)$.

Finally, it is straightforward to notice that if $G_{\pi_1} + G_{\pi_2}^{-1}$ has no unbalanced node, then $|Q_{\pi_1}| = |Q_{\pi_2}|$ and it is equal to the number of non-isolated nodes in this multigraph. \square

These propositions allow us to compute $TD(T_1, T_2)$, for $T_1, T_2 \in \mathcal{T}_n$, in time linear on n using the procedure given in pseudocode in Algorithm 2.

Remark 3. If T_1 and T_2 are two phylogenetic trees with different sets of labels, then we can compute their transposition distance by first restricting them to the sets of leaves with common labels, and then relabeling consecutively these common labels, starting with 1. Since we do not allow outdegree 1 nodes, when we restrict a phylogenetic tree to a subset of its set of taxa we contract edges to remove outdegree 1 nodes.

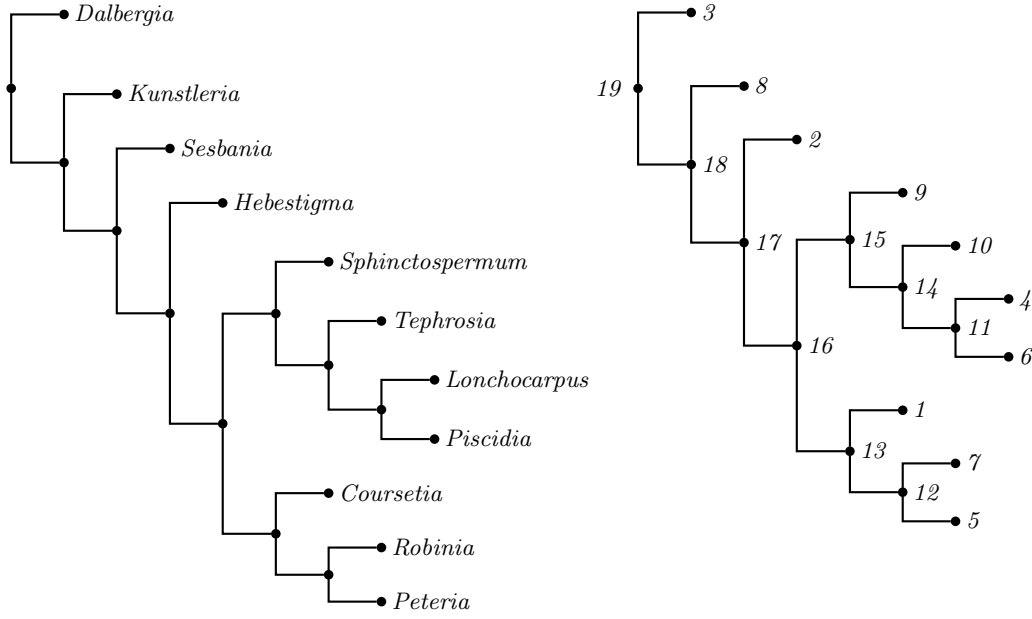


Fig. 4. A phylogenetic tree and the botom-up ordering of its restriction to the taxa of the tree in Fig. 1.

Example 5. Let T_1 be the phylogenetic tree in Example 1 and let T_2 be the lower phylogenetic tree displayed in Fig. 4, which represents the bottom-up ordering (with its taxa sorted alphabetically) of the tree T270c2x3x96c12c57c27 in TreeBASE after removing the outer taxon *Dalbergia* (and the elementary root created in this way), which does not appear in T_1 . Its matching permutation is

$$\pi(T_2) = (4, 6)(7, 5)(1, 12)(10, 11)(9, 14)(13, 15)(2, 16)(8, 17)(3, 18).$$

Since

$$\pi(T_1) = (1, 5, 7, 9)(4, 6, 10)(2, 11)(8, 13)(3, 12, 14),$$

(see Example 3), the multigraph $G_{\pi(T_1)} + G_{\pi(T_2)}^{-1}$ has nodes $\{1, \dots, 18\}$, red arcs $(1, 5)$, $(5, 7)$, $(7, 9)$, $(9, 1)$, $(4, 6)$, $(6, 10)$, $(10, 4)$, $(2, 11)$, $(11, 2)$, $(8, 13)$, $(13, 8)$, $(3, 12)$, $(12, 14)$, and $(14, 3)$, and blue arcs $(4, 6)$, $(6, 4)$, $(7, 5)$, $(5, 7)$, $(1, 12)$, $(12, 1)$, $(10, 11)$, $(11, 10)$, $(9, 14)$, $(14, 9)$, $(13, 15)$, $(15, 13)$, $(2, 16)$, $(16, 2)$, $(8, 17)$, $(17, 8)$, $(3, 18)$, and $(18, 3)$.

To compute $TD(T_1, T_2)$, we start with $d = 0$ and $N = 18$.

1. At the beginning, 15, 16, 17 and 18 are unbalanced. Then, we remove the pairs of blue arcs $\{(13, 15), (15, 13)\}$, $\{(2, 16), (16, 2)\}$, $\{(8, 17), (17, 8)\}$, and $\{(3, 18), (18, 3)\}$ and we set $d = 4$ and $N = 14$.

2. In this way, the nodes 2, 3, 8, 13 become unbalanced. Then, we remove the pairs of red arcs $\{(2, 11), (11, 2)\}$, $\{(8, 13), (13, 8)\}$ and we replace the pair of red arcs $(14, 3), (3, 12)$ by a new red arc $(14, 12)$ and we set $d = 7$ and $N = 10$.
3. Now, 11 has become unbalanced. Then, we remove the pair of blue arcs $\{(10, 11), (11, 10)\}$ and we set $d = 8$ and $N = 9$.
4. Now, 10 has become unbalanced. Then, we replace the pair of red arcs $(6, 10), (10, 4)$ by a new red arc $(6, 4)$ and we set $d = 9$ and $N = 8$.
5. At this moment, there does not remain any unbalanced node: the resulting multigraph has 5 alternating cycles (a cycle $(1, 5, 7, 9, 14, 12, 1)$, a cycle $(1, 12, 14, 9, 1)$, a cycle $(5, 7, 5)$, and two cycles $(4, 6, 4)$). Then, we have

$$TD(T_1, T_2) = \frac{1}{2}(d + N - 5) = 6.$$

In the Introduction we mentioned that the transposition distance defined in this paper generalizes the transposition distance for fully resolved trees. This will be a direct consequence of the following result.

Proposition 5. *For every pair of binary phylogenetic trees $T_1, T_2 \in \mathcal{T}_n$, let $G = (V, E)$ be the undirected multigraph with $V = \{1, \dots, 2n - 2\}$ and $E = M(T_1) \sqcup M(T_2)$, and let C be the set of connected components of G . Then, $TD(T_1, T_2) = n - 1 - |C|$.*

Proof. Let G_1 and G_2 denote the directed graphs associated to $\pi(T_1)$ and $\pi(T_2)^{-1}$. Since T_1 and T_2 are binary, in $G_1 + G_2$ for every blue or red arc (i, j) there is the inverse arc (j, i) of the same color, and the graph G in the statement is the undirected graph obtained by replacing each pair of arcs of the same color $\{(i, j), (j, i)\}$ by the undirected edge $\{i, j\}$, which we shall understand colored with the same color as the original pair.

Since $T_1, T_2 \in \mathcal{T}_n$ are binary, and therefore they have $2n - 1$ nodes, no one of the $2n - 2$ nodes of $G_1 + G_2$ is unbalanced or isolated. Then, by Proposition 4,

$$TD(T_1, T_2) = \frac{1}{2}((2n - 2) - A(G_1, G_2)) = n - 1 - \frac{1}{2}A(G_1, G_2).$$

Moreover, G is 2-regular, and therefore, every connected component in G is an alternating cycle, which contains exactly two alternating cycles of $G_1 + G_2$. Therefore $A(G_1, G_2) = 2|C|$. Combining this equality with the expression for $TD(T_1, T_2)$ given by Proposition 4, we obtain the expression in the statement. \square

In [17], the *transposition distance* between two binary phylogenetic trees T_1 and T_2 was defined as the least number of transpositions necessary to transform $M(T_1)$ into $M(T_2)$: in this context, a *transposition* means a replacement of a pair of 2-elements sets $\{i, j\}, \{k, l\}$ by a new pair $\{i, k\}, \{j, l\}$. Theorem 1 in *loc. cit.* and the last proposition entail that, for binary phylogenetic trees, our transposition distance and the transposition distance defined in [17] are the same.

4 Results

We have implemented in Perl the algorithms for the transposition distance between phylogenetic trees, using the BioPerl collection of Perl modules for computational biology [15]. The software is available in source code form for research use to educational institutions, non-profit research institutes, government research laboratories, and individuals, for non-exclusive use, without the right of the licensee to further redistribute the source code. The software is also provided for free public use on a web server, at the address <http://www.lsi.upc.edu/~valiente/>

Using this implementation, we have performed a systematic study of the TreeBASE [7] phylogenetic database, the main repository of published phylogenetic analyses, which currently contains 2,592 phylogenies with 36,593 taxa among them. Previous studies have revealed that TreeBASE constitutes a scale-free network [10].

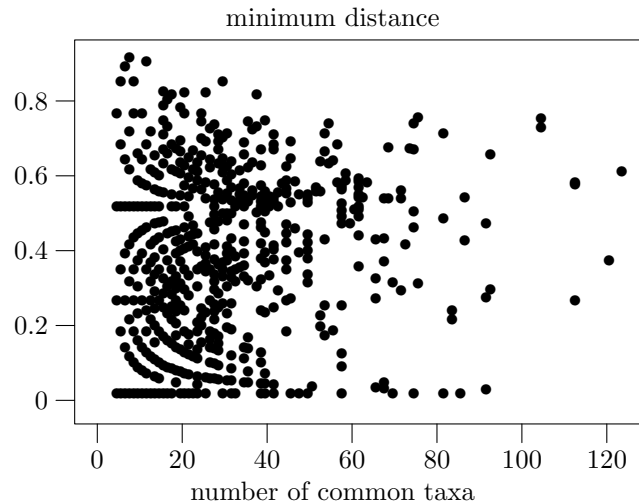


Fig. 5. Similarity of phylogenetic trees in TreeBASE based on the transposition distance. Each bullet represents the distance between a phylogenetic tree and the most similar phylogenetic tree in TreeBASE (other than itself) with at least three common taxa.

In order to assess the usefulness of the new distance measure in practice, we have computed the transposition distance for each of the $2,592 \cdot 2,591/2 = 3,357,936$ pairs of phylogenetic trees in TreeBASE. Then, for each phylogenetic tree, we have recovered the most similar phylogenetic tree in TreeBASE (other than itself) with at least three taxa in common. The results, summarized in Fig. 5, show that the transposition distance allows for a good recall of similar phylogenetic trees.

Acknowledgements

This work has been partially supported by the Spanish DGES project BFM2003-00771 ALBIOM, the Spanish CICYT project TIN 2004-07925-C03-01 GRAMMARS, and the UE project INTAS IT 04-77-7178.

References

1. B. L. Allen, Mike A. Steel. "Subtree transfer operations and their induced metrics on evolutionary trees." *Ann. Combin.*, 5 (2001), 1–13.
2. The American Institute of Mathematics. "Geometric models of biological phenomena." <http://www.aimath.org/WWN/geombio/geombio.pdf> (2003).
3. J. Bluis, D.-G. Shin. "Nodal distance algorithm: Calculating a phylogenetic tree comparison metric." In *Proc. 3rd IEEE Symposium on BioInformatics and BioEngineering* (2003) 87–94.
4. P. W. Diaconis, S. P. Holmes. "Matchings and phylogenetic trees." *Proc. Natl. Acad. Sci. USA*, 95 (1998), 14600–14602.
5. G. Estabrook, F. McMorris, C. Meacham. "Comparison of undirected phylogenetic trees based on subtrees of four evolutionary units." *Syst. Zool.*, 34 (1985), 193–200.
6. J. Handl, J. Knowles, D. B. Kell. "Computational cluster validation in post-genomic data analysis." *Bioinformatics*, 21 (2005), 3201–3212.
7. V. Morell. "TreeBASE: The roots of phylogeny." *Science*, 273 (1996), 569–570. <http://www.treebase.org>
8. R. D. M. Page. "Phyloinformatics: Towards a phylogenetic database." In *Data Mining in Bioinformatics*, chapter 10 (Springer-Verlag, 2005), 219–241.
9. D. Penny, M. D. Hendy. "The use of tree comparison metrics." *Syst. Zool.*, 34 (1985), 75–82, .
10. W. H. Piel, M. J. Sanderson, M. J. Donoghue. "The small-world dynamics of tree networks and data mining in phyloinformatics." *Bioinformatics*, 19 (2003), 1162–1168.
11. C. Reidys, P. F. Stadler. "Bio-molecular shapes and algebraic structures." *Computers & Chemistry*, 20 (1996), 85–94.
12. D. F. Robinson, L. R. Foulds. "Comparison of weighted labelled trees." In *Proc. 6th Australian Conf. Combinatorial Mathematics*, Lecture Notes in Mathematics 748 (1979), 119–126.
13. D. F. Robinson, L. R. Foulds. "Comparison of phylogenetic trees." *Math. Biosci.*, 53 (1981), 131–147.
14. F. Rosselló. "Reidys's and Stadler's metrics for RNA secondary structures." *Mathematical and Computer Modelling*, 40 (2004), 771–776. <http://es.arxiv.org/abs/math.GM/0305222>
15. J. E. Stajich, D. Block, K. Boulez, S. E. Brenner, S. A. Chervitz, C. Dagdigian, G. Fuellen, J. G. Gilbert, I. Korf, H. Lapp, H. Lehvaslaiho, C. Matsalla, C. J. Mungall, B. I. Osborne, M. R. Pocock, P. Schattner, M. Senger, L. D. Stein, E. Stupka, M. D. Wilkinson, E. Birney. "The BioPerl toolkit: Perl modules for the life sciences." *Genome Research*, 12 (2002), 1611–1618. <http://www.bioperl.org>
16. R. P. Stanley. *Enumerative combinatorics, Vol. 2. Cambridge Studies in Advanced Mathematics* vol. 62, Cambridge Univ. Press (1998).
17. G. Valiente. "A fast algorithmic technique for comparing large phylogenetic trees." In *Proc. 12th Int. Symp. String Processing and Information Retrieval*, Lecture Notes in Mathematics 3772 (2005), 370–375.
18. Gabriel Valiente. *Algorithms on Trees and Graphs*. Springer-Verlag (2002).
19. M. S. Waterman, T. F. Smith. "On the similarity of dendograms." *J. Theor. Biol.*, 73 (1978) 789–800.